

Selecting optimal minimum spanning trees that share a topological correspondence with phylogenetic trees.

Prabhav Kalaghatgi

Max Planck Institute for Informatics
Saarbrücken
prabhavk@mpi-inf.mpg.de

Thomas Lengauer

Max Planck Institute for Informatics
Saarbrücken
lengauer@mpi-inf.mpg.de

Abstract

Choi *et al.* (2011) introduced a minimum spanning tree (MST)-based method called CLGrouping, for constructing tree-structured probabilistic graphical models, a statistical framework that is commonly used for inferring phylogenetic trees. While CLGrouping works correctly if there is a unique MST, we observe an indeterminacy in the method in the case that there are multiple MSTs. In this work we remove this indeterminacy by introducing so-called vertex-ranked MSTs. We note that the effectiveness of CLGrouping is inversely related to the number of leaves in the MST. This motivates the problem of finding a vertex-ranked MST with the minimum number of leaves (MLVRMST). We provide a polynomial time algorithm for the MLVRMST problem, and prove its correctness for graphs whose edges are weighted with tree-additive distances.

1 Introduction

Phylogenetic trees are commonly modeled as tree-structured probabilistic graphical models with two types of vertices: labeled vertices that represent observed taxa, and hidden vertices that represent unobserved ancestors. The length of each edge in a phylogenetic tree quantifies evolutionary distance. If the set of taxa under consideration contain ancestor-descendant pairs, then the phylogenetic tree has labeled internal vertices, and is called a *generally labeled tree* (Kalaghatgi *et al.*, 2016). The data that is used to infer the topology and edge lengths is usually available in the form of gene or protein sequences.

Popular distance-based methods like neighbor joining (NJ; Saitou and Nei (1987)) and BIONJ (Gascuel, 1997) construct phylogenetic trees from estimates of the evolutionary distance between each pair of taxa. Choi *et al.* (2011) introduced a distance-based method called Chow-Liu grouping (CLGrouping). Choi *et al.* (2011) argue that CLGrouping is more accurate than NJ at reconstructing phylogenetic trees with large diameter. The diameter of tree is the number of edges in the longest path of the tree.

CLGrouping operates in two phases. The first phase constructs a distance graph G which is a complete graph over the labeled vertices where each edge is weighted with the distance between each pair of labeled vertices. Subsequently a minimum spanning tree (MST) of G is constructed. In the second phase, for each internal vertex v_i of the MST, the vertex set V_i consisting of v_i and its neighbors is constructed. Subsequently a generally labeled tree T_i over V_i is inferred using a distance-based tree construction method like NJ. The subtree in the MST that is induced by V_i is replaced with T_i .

Distances are said to be additive in a tree T if the distance between each pair of vertices u and v is equal to the sum of lengths of edges that lie on the path in T between u and v . Consider the set of all phylogenetic trees \mathcal{T} such that the edge length of each edge in each tree in \mathcal{T} is strictly greater than zero. A distance-based tree reconstruction method is said to be consistent if for each $\{D, T | T \in \mathcal{T}\}$ such that D is additive in T , the tree that is reconstructed using D is identical to T . Please note the following well-known result regarding the correspondence between trees and additive distances. Considering all trees in \mathcal{T} , if D is additive in a tree T then T is unique (Buneman, 1971).

We show that if G has multiple MSTs then CLGrouping is not necessarily consistent. We show that there always exists an MST M such that CLGrouping returns the correct tree when M is used in the second phase of CLGrouping. We show that M can be constructed by assigning ranks to the vertices in G , and by modifying standard MST construction algorithms such that edges are compared on the basis of both edge weight and ranks of the incident vertices. The MSTs that are constructed in this manner are called vertex-ranked MSTs.

Given a distance graph, there may be multiple vertex-ranked MSTs with vastly different number of leaves. Huang *et al.* (2014) showed that CLGrouping affords a high degree of parallelism, because, phylogenetic tree reconstruction for each vertex group can be performed independently. With respect to parallelism, we define an optimal vertex-ranked MST for CLGrouping to be a vertex-ranked MST with the maximum number of vertex groups, and equivalently, the minimum number of leaves.

We developed an $O(n^2 \log n)$ time algorithm Algo. 1 that takes as input a distance graph and outputs a vertex-ranked MST with the minimum number of leaves (MLVRMST). The proof of correctness of Algo. 1 assumes that the edges in the distance graph are weighted with tree-additive distances.

2 Terminology

A phylogenetic tree is an undirected edge-weighted acyclic graph with two types of vertices: labeled vertices that represent observed taxa, and hidden vertices that represent unobserved taxa. Information, e.g., in the form of genomic sequences, is only present at labeled vertices. We refer to the edge weights of a phylogenetic tree as edge lengths. The length of an edge quantifies the estimated evolutionary distance between the sequences corresponding to the respective incident vertices. All edge lengths are strictly positive. Trees are leaf-labeled if all the labeled vertices are leaves. Leaf-labeled phylogenetic trees are the most commonly used models of evolutionary relationships. Generally labeled trees are phylogenetic trees whose internal vertices may be labeled, and are appropriate when ancestor-descendant relationships may be present in the sampled taxa (Kalaghatgi *et al.*, 2016).

Each edge in a phylogenetic tree partitions the set of all labeled vertices into two disjoint sets which are referred to as the split of the edge. The two disjoint sets are called to the sides of the split.

A phylogenetic tree can be rooted by adding a hidden vertex (the root) to the tree, removing an edge e in the tree, and adding edges between the root and the vertices that were previously incident to e . Edge lengths for the newly added edges must be positive numbers and must sum up to the edge length of the previously removed edge. Rooting a tree constructs a directed acyclic graph in which each edge is directed away from the root.

A leaf-labeled phylogenetic tree is clock-like if the tree can be rooted in such a way that all leaves are equidistant from the root. Among all leaf-labeled phylogenetic trees, maximally balanced trees and caterpillar trees have the smallest and largest diameter, respectively, where the diameter of a tree is defined as the number of edges along the longest path in the tree.

The distance graph G of a phylogenetic tree T is the edge-weighted complete graph whose vertices are the labeled vertices of T . The weight of each edge in G is equal to the length of the path in T that connects the corresponding vertices that are incident to the edge. A minimum spanning tree (MST) of an edge-weighted graph is a tree that spans all the vertices of the graph, and has the minimum sum of edge weights.

3 Chow-Liu grouping

Choi *et al.* (2011) introduced the procedure Chow-Liu grouping (CLGrouping) for the efficient reconstruction of phylogenetic trees from estimates of evolutionary distances. If the input distances are additive in the phylogenetic tree T then the authors claim that CLGrouping correctly reconstructs T .

CLGrouping consists of two stages. In the first stage, an MST M of G is constructed. In the second stage, for each internal vertex v , a vertex group $Nb(v)$ is defined as follows: $Nb(v)$ is the set containing v and all the vertices in M that are adjacent to v . For each vertex group, a phylogenetic tree T_v is constructed using

distances between vertices in $Nb(v)$. Subsequently, the graph in M that is induced by $Nb(v)$ is replaced by T_v (see Fig. 1e for an illustration). T_v may contain hidden vertices which may now be in the neighborhood of an internal vertex w that has not been visited as yet. If this the case, then we need an estimate of the distance between the newly introduced hidden vertices and vertices in $Nb(w)$. Let h_v be the hidden vertex that was introduced when processing the internal vertex v . The distance from h_v to a vertex $k \in Nb(w)$ is estimated using the following formula, $d_{h_v k} = d_{vk} - d_{vh_v}$.

The order in which the internal vertices are visited is not specified by the authors and does not seem to be important. CLGrouping terminates once all the internal vertices of M have been visited.

This procedure is called Chow-Liu grouping because the MSTs that are constructed using additive distances are equivalent to Chow-Liu trees (Chow and Liu, 1968), for certain probability distributions. Please read Choi *et al.* (2011) for further detail.

4 Indeterminacy of CLGrouping

CLGrouping is not necessarily consistent if there are multiple MSTs. We demonstrate this with the phylogenetic tree T shown in Fig. 1a. For the corresponding distance graph G of T (see Fig. 1b), two MSTs of G , M_1 and M_2 are shown in Fig. 1c and Fig. 1d, respectively. The intermediate steps, and the final result of applying CLGrouping to M_1 and M_2 are shown in Fig. 1e and Fig. 1f, respectively. CLGrouping reconstructs the original phylogenetic tree if it is applied to M_1 but not if it is applied to M_2 .

The notion of a surrogate vertex is central to proving the correctness of CLGrouping. The surrogate vertex of a hidden vertex is the closest labeled vertex, w.r.t. distances defined on the phylogenetic tree. CLGrouping will reconstruct the correct phylogenetic tree only if the MST can be constructed by contracting all the edges along the path between each hidden vertex and its surrogate vertex. Since the procedure that constructs the MST is not aware of the true phylogenetic tree, the surrogate vertex of each hidden vertex must be selected implicitly. In the example shown earlier, M_1 can be constructed by contracting the edges (h_1, l_1) , and (h_2, l_3) . Clearly there is no selection of surrogate vertices such that M_2 can be constructed by contracting the path between each hidden vertex and the corresponding surrogate vertex.

If there are multiple labeled vertices each of which is closest to a hidden vertex then Choi *et al.* (2011) assume that the corresponding surrogate vertex is implicitly selected using the following tie-breaking rule.

Let the surrogate vertex set $\mathbf{Sg}(h)$ of a vertex h be the set of all labeled vertices that are closest to h . If l_1 and l_2 belong to both $\mathbf{Sg}(h_1)$ and $\mathbf{Sg}(h_2)$, then the same labeled vertex (either l_1 or l_2) is selected as the surrogate vertex of both h_1 and h_2 . This rule for selecting surrogate vertices cannot be consistently applied across all hidden vertices. We demonstrate this with an example. For the tree shown in Fig. 2 we have $\mathbf{Sg}(h_1) = \{l_1, l_2\}$, $\mathbf{Sg}(h_2) = \{l_4, l_5\}$, and $\mathbf{Sg}(h_3) = \{l_1, l_2, l_3, l_4, l_5\}$. It is clear that there is no selection of surrogate vertices that satisfies the tie-breaking rule.

5 Ensuring the consistency of CLGrouping

In order to construct an MST that is guaranteed to have the desired topological correspondence with the phylogenetic tree, we propose the following tie-breaking rule for selecting the surrogate vertex. Let there be a total order over the set of all labeled vertices. Let $\mathcal{R}(l)$ be the rank of vertex l that is given by the order. We define the surrogate vertex $\mathbf{Sg}(h)$ of h to be the highest ranked labeled vertex among the set of labeled vertices that are closest to h . That is,

Definition 1.

$$\mathbf{Sg}(h) = \min_{l \in \mathbf{Sg}(h)} \mathcal{R}(l), \text{ where,}$$

$$\mathbf{Sg}(h) = \min_{l \in \mathcal{L}(T)} d_{lh}.$$

The inverse surrogate set $\mathbf{Sg}^{-1}(l)$ is the set of all hidden vertices whose surrogate vertex is l .

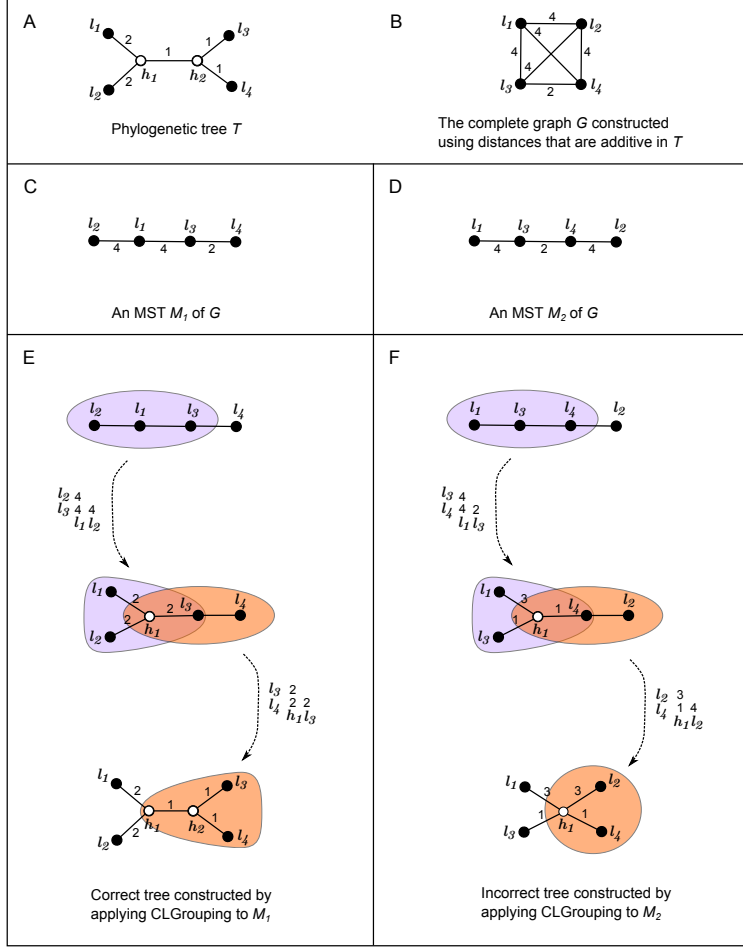


Figure 1: The example used to demonstrate that CLGrouping may not reconstruct the correct tree if there are multiple MSTs. The phylogenetic tree T that is used in this example is shown in panel a. The distance graph G of T is shown in panel b. Two MSTs of G , M_1 and M_2 , respectively, are shown in panels c and d. Panels e and f show the intermediate steps and the final result of applying CLGrouping to M_1 and M_2 respectively. CLGrouping reconstructs the original phylogenetic tree if it is applied to M_1 , but not if it is applied to M_2 .

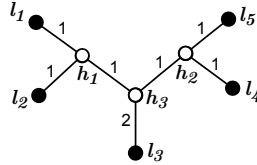


Figure 2: The phylogenetic tree that is used to demonstrate that the tie-breaking rule as defined by Choi *et al.* (2011) cannot be applied in general.

In order to ensure that the surrogate vertices are selected on the basis of both distance from the corresponding hidden vertex and vertex rank, it is necessary that information pertaining to vertex rank is used when selecting the edges of the MST. We use Kruskal's algorithm (Kruskal, 1956) for constructing the desired MST. Since Kruskal's algorithm takes as input a set of edges sorted w.r.t. edge weight, we modify the

input by sorting edges with respect to edge weight and vertex rank as follows. It is easy to modify other algorithms for constructing MSTs in such a way that vertex rank is taken into account.

Definition 2. We define below, what is meant by sorting edges on the basis of edge weight and vertex rank. Given a edge set E , and a ranking \mathcal{R} over vertices in E , let $d(u, v)$ be the weight of the edge $\{u, v\}$, and let $\mathcal{R}(u)$ be the rank of the vertex u . Let the relative position of each pair of edges in the list of sorted edges be defined using the total order $<$. That is to say, for each pair of edges $\{a, b\}$ and $\{c, d\}$,

$\{a, b\} < \{c, d\}$, if and only if

(i) $d(a, b) < d(c, d)$, or if

(ii) $d(a, b) = d(c, d)$ and $\min\{\mathcal{R}(a), \mathcal{R}(b)\} < \min\{\mathcal{R}(c), \mathcal{R}(d)\}$, or if

(iii) $d(a, b) = d(c, d)$ and $\min\{\mathcal{R}(a), \mathcal{R}(b)\} = \min\{\mathcal{R}(c), \mathcal{R}(d)\}$ and $\max\{\mathcal{R}(a), \mathcal{R}(b)\} < \max\{\mathcal{R}(c), \mathcal{R}(d)\}$.

The MST that is constructed by applying Kruskal's algorithm to the edges that are ordered with respect to weight and vertex rank is called a vertex-ranked MST (VRMST).

Now, we will prove Lemma 1, which is used to prove the correctness of CLGrouping.

Lemma 1. *Adapted from parts (i) and (ii) of Lemma 8 in Choi et al. (2011). Given a phylogenetic tree T and a ranking \mathcal{R} over the labeled vertices in T , let G be the distance graph that corresponds to $T = (V_T, E_T)$ and let E_{\leq} be the list of edges of G sorted with respect to edge weight and vertex rank, as defined in Definition 2. Let $M = (V_M, E_M)$ be the VRMST that is constructed by applying Kruskal's algorithm to E_{\leq} . The surrogate vertex of each hidden vertex is defined with respect to distance and vertex rank as given in Definition 1. M is related to T as follows.*

(i) *If $j \in V_M$ and $h \in Sg^{-1}(j)$ s.t. $h \neq j$, then every vertex in the path in T that connects j and h belongs to the inverse surrogate set $Sg^{-1}(j)$.*

(ii) *For any two vertices that are adjacent in T , their surrogate vertices, if distinct, are adjacent in M , i.e., for all $i, j \in V_T$ with $Sg(i) \neq Sg(j)$,*

$$\{i, j\} \in E_T \Rightarrow \{Sg(i), Sg(j)\} \in E_M$$

Proof 1. First we will prove Lemma 1 part (i) by contradiction.

Assume that there is a vertex u on the path between h and j , such that $Sg(u) = k \neq j$. We have $d_{uk} \leq d_{uj}$ (equality holds only if $\mathcal{R}(k) < \mathcal{R}(j)$). Similarly, since $Sg(h) = j$, we have $d_{hj} \leq d_{hk}$ (equality holds only if $\mathcal{R}(j) < \mathcal{R}(k)$). We consider all eight positions of k w.r.t. h, u , and j (see Fig. 3).

For case 1 we have

$$\begin{aligned} d_{hj} &\leq d_{hk} \text{ (since } Sg(h) = j) \\ \Leftrightarrow d_{hl} + d_{lu} + d_{uj} &\leq d_{hl} + d_{hk} \\ \Leftrightarrow d_{lu} + d_{uj} &\leq d_{lk} \\ \Leftrightarrow d_{uj} &< d_{ul} + d_{lk} \\ \Leftrightarrow d_{uj} &< d_{uk} \text{ (contradiction since } Sg(u) = k). \end{aligned}$$

For case 2 we have

$$\begin{aligned} d_{hj} &\leq d_{hk} \text{ (equality holds only if } \mathcal{R}(j) < \mathcal{R}(k)) \\ \Leftrightarrow d_{hu} + d_{ul} + d_{lj} &\leq d_{hu} + d_{ul} + d_{lk} \\ \Leftrightarrow d_{ul} + d_{lj} &\leq d_{ul} + d_{lk} \\ \Leftrightarrow d_{uj} &\leq d_{uk} \text{ (contradiction since } Sg(u) = k). \end{aligned}$$

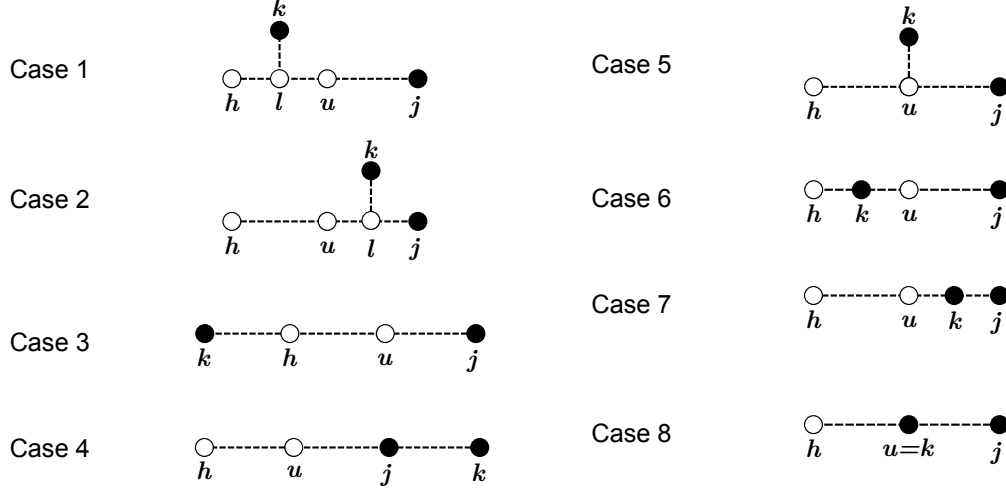


Figure 3: The cases that were considered in the proof of Lemma 1 part (ii). For some phylogenetic tree T let j be a labeled vertex and let h be a hidden vertex in the inverse surrogate set of j . u is a vertex in the path between h and j . Each case specifies one of the eight possible positions of a labeled vertex k w.r.t h, u , and j . Hidden vertices are represented with white circles and labeled vertices are represented with black circles. Each dashed line represents a path between the two vertices at its end points.

For case 3 we have

$$\begin{aligned}
d_{hj} &\leq d_{hk} \\
\Leftrightarrow d_{hu} + d_{uj} &\leq d_{hk} \\
\Leftrightarrow d_{uj} &< d_{hk} + d_{hu} \\
\Leftrightarrow d_{uj} &< d_{uk} \text{ (contradiction since } \text{Sg}(u) = k\text{)}.
\end{aligned}$$

For case 4 we have

$$\begin{aligned}
d_{uk} &= d_{uj} + d_{jk} \text{ (see Fig. 3 case 4)} \\
\Leftrightarrow d_{uk} &> d_{uj} \text{ (contradiction since } \text{Sg}(u) = k\text{)}.
\end{aligned}$$

For case 5 we have

$$\begin{aligned}
d_{hj} &\leq d_{hk} \text{ (equality holds only if } \mathcal{R}(j) < \mathcal{R}(k)\text{)} \\
\Leftrightarrow d_{hu} + d_{uj} &\leq d_{hu} + d_{uk} \\
\Leftrightarrow d_{uj} &\leq d_{uk} \text{ (contradiction since } \text{Sg}(u) = k\text{)}.
\end{aligned}$$

For cases 6,7, and 8, we have

$$\begin{aligned}
d_{hj} &\leq d_{hk} \\
\Leftrightarrow d_{hk} + d_{kj} &\leq d_{hk} \\
\Leftrightarrow d_{hk} &< d_{hk} \text{ (contradiction)}.
\end{aligned}$$

Now we will prove part (ii) of Lemma 1. Consider the edge $\{i, j\}$ in E_T such that $\text{Sg}(i) \neq \text{Sg}(j)$. Let V_i and V_j be the sides of the split that is induced by the edge $\{i, j\}$, such that V_i and V_j contain i and j , respectively. Let L_i and L_j be sets of labeled vertices that are defined as $V_i \cap V_M$ and $V_j \cap V_M$ respectively.

From part (i) of Lemma 1 we know that $\text{Sg}(i) \in L_i$ and $\text{Sg}(j) \in L_j$. Additionally, for any $k \in L_i \setminus \{\text{Sg}(i)\}$ and $l \in L_j \setminus \{\text{Sg}(j)\}$, from the definition of surrogate vertex it follows that

$$\begin{aligned} d_{ki} &\geq d_{\text{Sg}(i)i} \text{ (equality holds only if } \mathcal{R}(\text{Sg}(i)) < \mathcal{R}(k)) \\ d_{lj} &\geq d_{\text{Sg}(j)j} \text{ (equality holds only if } \mathcal{R}(\text{Sg}(j)) < \mathcal{R}(l)) \\ d_{kj} &= d_{ki} + d_{ij} + d_{lj} \\ &\geq d_{\text{Sg}(i)i} + d_{ij} + d_{\text{Sg}(j)j} \\ &= d_{\text{Sg}(i)\text{Sg}(j)}. \end{aligned}$$

It is clear that

$$\min\{\mathcal{R}(k), \mathcal{R}(l)\} > \min\{\mathcal{R}(\text{Sg}(i)), \mathcal{R}(\text{Sg}(j))\}, \quad (1)$$

and that

$$d_{kl} \geq d_{\text{Sg}(i)\text{Sg}(j)}. \quad (2)$$

The cut property of MSTs states that given a graph $G = (V, E)$ for each pair V_1, V_2 of disjoint sets such that $V_1 \cup V_2 = V$, each MST of G contains one of the smallest edges (w.r.t. edge weight) which have one end-point in V_1 and the other end-point in V_2 .

Note that the vertex-ranked MST M is constructed using edges that are sorted w.r.t. edge weight and the vertex rank \mathcal{R} . From equations (1) and (2) it is clear that among all edges with one end point in L_i and the other end-point in L_j , the edge $\{\text{Sg}(i), \text{Sg}(j)\}$ is the smallest edge w.r.t edge weight and vertex rank (see Definition 2). Since L_i and L_j are disjoint sets and $L_i \cup L_j = V_M$, it follows that $\{\text{Sg}(i), \text{Sg}(j)\} \in E_M$. \square

CLGrouping can be shown to be correct using Lemma 1 and the rest of the proof that was provided by Choi *et al.* (2011). Thus if the distances are additive in the model tree, CLGrouping will provably reconstruct the model tree provided that the MST that is used by CLGrouping is a vertex-ranked MST (VRMST).

The authors of CLGrouping provide a matlab implementation of their algorithm. Their implementation reconstructs the model tree even if there are multiple MSTs in the underlying distance graph. The authors' implementation takes as input a distance matrix which has the following property: the row index, and the column index of each labeled vertex is equal. The MST that is constructed in the authors implementation is a vertex-ranked MST, with the rank of each vertex being equal to the corresponding row index of the labeled vertex. We implemented their algorithm in python with no particular order over the input distances and were surprised to find out that the reconstructed tree differed from the model tree, even if the input distances were additive in the model tree.

Depending on the phylogenetic tree, there may be multiple corresponding vertex-ranked MSTs with vastly different numbers of leaves. In the next section we discuss the impact of the number of leaves in a vertex-ranked MST, on the efficiency of parallel implementations of CLGrouping.

6 Relating the number of leaves in a VRMST to the optimality of the VRMST in the context of CLGrouping

In the context of parallel programming, Huang *et al.* (2014) showed that it is possible to parallelize CLGrouping by independently constructing phylogenetic trees for each vertex group, and later combining them in order to construct the full phylogenetic tree.

In order to relate the balancedness of a phylogenetic tree to the number of leaves in a corresponding vertex-ranked MST, we consider clock-like caterpillar trees and maximally balanced trees such that each hidden vertex of each tree has degree three.

Consider the case in which the phylogenetic tree is a caterpillar tree (least balanced). There exists a corresponding VRMST which has a star topology that can be constructed by contracting edges between each

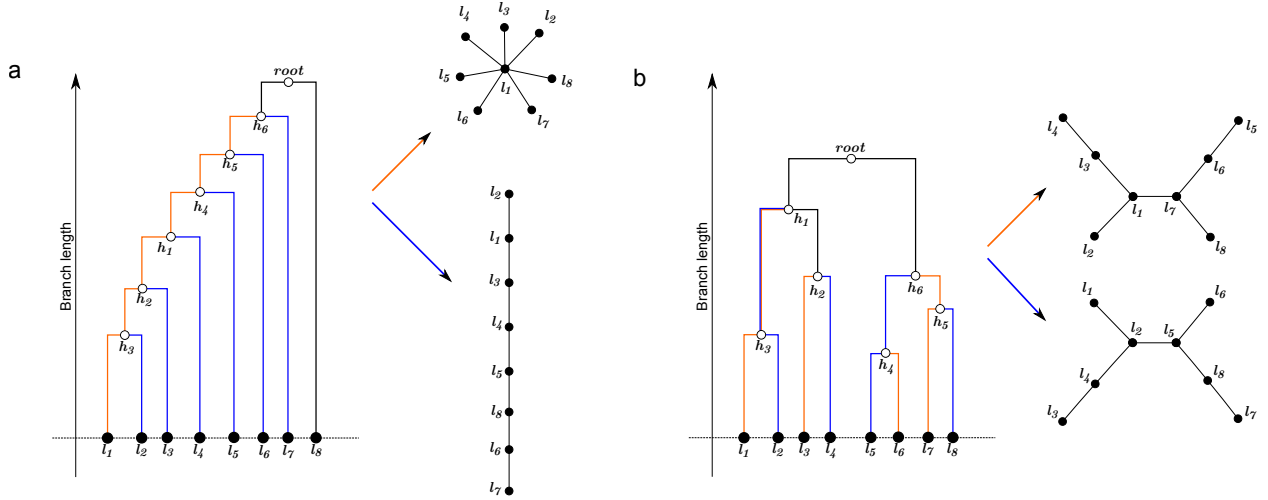


Figure 4: Both panels show clock-like phylogenetic trees and VRMSTs with the maximum and the minimum number of leaves, that are constructed by contracting corresponding edges that are highlighted in orange and blue, respectively. The difference between the maximum and the minimum number of leaves in VRMSTs is largest for the caterpillar tree shown in panel a, and smallest for the maximally balanced tree shown in panel b.

hidden vertex and one labeled vertex that is in the surrogate vertex set of each hidden vertex (see Fig 4a). A star-shaped VRMST has only one vertex group, comprising all the vertices in the VRMST, and does not afford any parallelism.

Instead, if the VRMST was to be constructed by contracting edges between each hidden vertex h and a labeled vertex that is incident to h , then the number of the vertex groups would be $n - 2$, where n is the number of vertices in the phylogenetic tree. The resulting VRMST would have the minimum number of leaves (two).

With respect to parallelism, an optimal vertex-ranked MST for CLGrouping is a vertex-ranked MST with the maximum number of vertex groups, and equivalently, the minimum number of leaves.

Consider a phylogenetic tree $T = (V_T, E_T)$ which is maximally balanced. It is clear that the set $\mathcal{L}(T)$ of labeled vertices of T can be partitioned into a disjoint set \mathcal{C} of vertex pairs such that for each vertex pair $\{u, v\} \in \mathcal{C}$, u and v are adjacent to the same hidden vertex $h \in V_T$. Given a vertex ranking \mathcal{R} , the surrogate vertex of h will be $\max_{l \in \{u, v\}} \mathcal{R}(l)$. Thus, independently of vertex ranking, the number of distinct surrogate vertices will be $\mathcal{L}(T)/2$. Each labeled vertex that is not selected as a surrogate vertex will be a leaf in the vertex-ranked MST. It follows that all corresponding VRMSTs of T will have $\mathcal{L}(T)/2$ leaves (see Fig 4b).

Whether or not the phylogenetic trees that are estimated from real data are clock-like depends on the set of taxa that are being studied. Genetic sequences that are sampled from closely related taxa have been estimated to undergo substitutions at a similar rate, resulting in clock-like phylogenetic trees (dos Reis *et al.*, 2016). In the context of evolution, trees are caterpillar-like if there is a strong selection; the longest path from the root represents the best-fit lineage.

In the next section we will present an algorithm for constructing a vertex-ranked MST with the minimum number of leaves.

7 Constructing a vertex-ranked MST with the minimum number of leaves

We aim to construct a vertex-ranked MST with the minimum number of leaves (MLVRMST) from a distance graph. An algorithm for constructing a MLVRMST is presented in subsection 7.3. In the following two subsections we will present two lemmas, which will be used for proving the correctness of the algorithm.

7.1 A common structure that is shared by all MSTs

In this section we will prove the existence of a laminar family \mathcal{F} over the vertex set of an edge-weighted graph G . A collection \mathcal{F} of subsets of a set S is a laminar family over S if, for any two intersecting sets in \mathcal{F} , one set contains the other. That is to say, for each pair S_1, S_2 in \mathcal{F} such that $|S_1| \leq |S_2|$, either $S_1 \cap S_2 = \emptyset$, or $S_1 \subset S_2$.

The vertex sets in \mathcal{F} define a structure that is common to each MST of G . Furthermore, \mathcal{F} can be used to obtain an upper bound on the degree of each vertex in a MST. The notion of a laminar family has been utilized previously by Ravi and Singh (2006), for designing an approximation algorithm for the minimum-degree MST

Lemma 2. *Given an edge-weighted graph $G = (V, E)$ with k distinct weight classes $W = \{w_1, w_2, \dots, w_k\}$, and an MST M of G , let F_i be the forest that is formed by removing all edges in G that are heavier than w_i . Let \mathcal{C}_i be the collection comprising the vertex set of each component of F_i . Consider the collection \mathcal{F} which is constructed as follows: $\mathcal{F} = \{\cup_{i=1}^k \mathcal{C}_i\} \cup V$. The following is true:*

- (i) \mathcal{F} is a laminar family over V
- (ii) Each vertex set in \mathcal{F} induces a connected subgraph in each MST of G

Proof. (i). Consider any two vertex sets S_1 and S_2 in \mathcal{F} . Let w_1 and w_2 be the weights of the heaviest edges in the subgraphs of M that are induced by S_1 and S_2 , respectively. Let F_1 and F_2 be the forests that are formed by removing all edges in M that are heavier than w_1 and w_2 , respectively. Let \mathcal{C}_1 and \mathcal{C}_2 be the collections comprising the vertex set of each component in F_1 and F_2 , respectively.

It is clear that $S_1 \in \mathcal{C}_1$ and $S_2 \in \mathcal{C}_2$. Consider the case where $w_1 = w_2$. Since $\mathcal{C}_1 = \mathcal{C}_2$, it follows that $S_1 \cap S_2 = \emptyset$. If $w_1 \neq w_2$, then without loss of generality, let $w_1 < w_2$. F_2 can be constructed by adding to F_1 all edges in M that are no heavier than w_2 . Each component in F_1 that is not in F_2 induces a connected subgraph in exactly one component of F_2 . If $S_1 \in \mathcal{C}_1 \cap \mathcal{C}_2$ then $S_1 \cap S_2 = \emptyset$. Otherwise, if $S_1 \in \mathcal{C}_1 \setminus \mathcal{C}_2$, then S_1 is a subset of exactly one set in \mathcal{C}_2 . This implies that either $S_1 \subset S_2$, or $S_1 \cap S_2 = \emptyset$. Thus \mathcal{F} is a laminar family over V .

(ii). Let S_i be the vertex set of a component in the subgraph G_i of G that is created by removing all edges in G_i that are heavier than w_i . It is clear that S_i induces a connected subgraph in each minimum spanning forest of G_i . For each minimum spanning forest there is a corresponding MST of G , such that the minimum spanning forest can be constructed by removing from the MST all the edges are heavier than w_i . It follows that S_i induces a connected subgraph in each MST of G . \square

7.2 Selecting surrogate vertices on the basis of maximum vertex degree

Lemma 3. *We are given a phylogenetic tree T , the corresponding distance graph $G = (V, E)$, and the laminar family \mathcal{F} of the distance graph. Let the subgraph $g = (V_g, E_g)$ of G contain all edges that are present in at least one MST of G . Let h be a hidden vertex in T such that there is a leaf l in $\mathbf{Sg}(h)$, and h is incident to l . Let S_i be a vertex set in \mathcal{F} and let w_i be the corresponding edge weight. Then the following holds:*

- (i) Let J_v be the set of all vertices that are incident to vertex v in g . Let S_v be the smallest sub-collection of \mathcal{F} that covers J_v but not v . Among all MSTs, the maximum vertex degree $\delta_{\max}(v)$ of v is $|S_v|$.
- (ii) $\delta_{\max}(l) \leq \delta_{\max}(v)$ for each vertex v in $\mathbf{Sg}(h)$.

Proof. (i). Let $J_v = \{j_1, j_2, \dots, j_k\}$ be the set of all vertices that are incident to v . Let M be some MST of G . Let $\mathcal{S}_v = \{S_1, S_2, \dots, S_m\}$ be the smallest sub-collection of \mathcal{F} that covers J_v and does not include v . Let \mathcal{S}_v contain a set S_i that covers multiple vertices in J . Let j_1 and j_2 be any two vertices in S_i . Let w_i be the heaviest weight on the path that joins j_1 and j_2 in M . The edges $\{v, j_1\}$ and $\{v, j_2\}$ are heavier than w_i . If they were not, then we would have $v \in S_i$. Since v, j_1 and j_2 are on a common cycle, each MST of G can only contain one of the two edges $\{v, j_1\}$, and $\{v, j_2\}$. It follows that for each set $S_i \in \mathcal{S}_v$, each MST can contain at most one edge which is incident to v and to a vertex in S_i . Thus the maximum number of edges that can be incident to v in any MST is the number of vertex sets in \mathcal{S}_v , i.e., $\delta_{\max}(v) = |\mathcal{S}_v|$.

(ii). Let J_l and J_v be the set of all vertices that are incident to l and v in g , respectively. Let $j \in J_l \setminus \mathbf{Sg}(h)$. The weight of the edge $\{j, l\} \in E_g$ is given by d_{jl} . $d_{jh} > d_{vh}$ since $j \notin \mathbf{Sg}(h)$. Thus $d_{lj} > d_{lv}$, and consequently $v \in J_l$. We have $d_{jl} = d_{jh} + d_{hl} = d_{jh} + d_{hv} = d_{jv}$. Consider the MST $M = (V_M, E_M)$ that contains the edges $\{l, v\}$ and $\{l, h\}$. Consider the spanning tree M' that is formed by removing $\{l, h\}$ from E_M and adding $\{v, h\}$. M' and M have the same sum of edge weights. Thus we also have $j \in J_v$. Consequently $J_l \subseteq J_v$. Let \mathcal{S}_l and \mathcal{S}_v be the smallest sub-collections of \mathcal{F} such that \mathcal{S}_l covers J_l but does not contain l , and \mathcal{S}_v covers J_v but does not contain v . \mathcal{S}_v covers both J_l and J_v since $J_l \subseteq J_v$. Thus $|\mathcal{S}_l| \leq |\mathcal{S}_v|$. From part (i), we know that $|\mathcal{S}_l| = \delta_{\max}(l)$ and $|\mathcal{S}_v| = \delta_{\max}(v)$. Thus $\delta_{\max}(l) \leq \delta_{\max}(v)$. \square

7.3 Constructing a minimum leaves vertex-ranked MST

We now give an overview of Algo. 1. Algo. 1 takes as input a distance graph $G = (V, E)$ and computes δ_{\max} for each vertex in V . Subsequently, a ranking \mathcal{R} over V is identified such that vertices with lower δ_{\max} are assigned higher ranks. The output of Algo. 1 is the vertex-ranked MST which is constructed using \mathcal{R} . If G is weighted with tree-additive distances then the output of Algo. 1 is a vertex-ranked MST with the minimum number of leaves (MLVRMST).

An example of a phylogenetic tree, a corresponding MLVRMST, and the output MST M of Algo. 1, is shown in Fig. 5. M is superimposed with the following: the laminar family \mathcal{F} , the subgraph g , and δ_{\max} for each vertex.

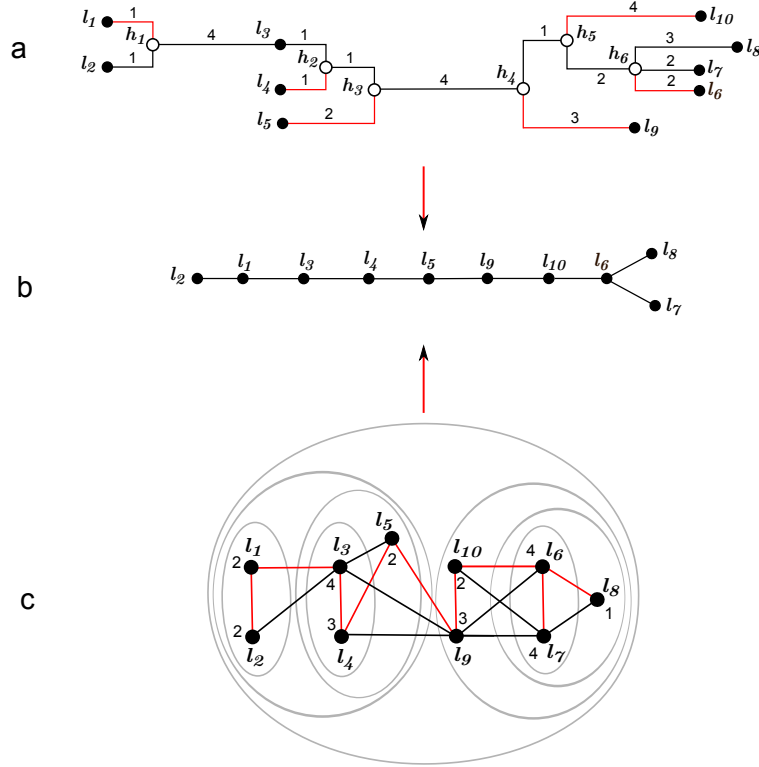


Figure 5: Panel a shows a generally labeled phylogenetic tree T . Algo. 1 was applied to the distance graph G of T . Panel b show the output M of Algo. 1 which is a MLVRMST of G . Panel c shows M (in red) superimposed with the laminar family \mathcal{F} , and the graph g which contains all edges of G that are present in at least one MST. Additionally each vertex has been labeled with the corresponding δ_{\max} .

Algorithm 1: MinLeavesVertexRankedMST of G

Input: $G = (V, E)$
Initialize:
For each vertex in V , create a Make-Set object;
Create empty arrays called E_w , W , E_{fixed} , $E_{flexible}$, and $E_{selected}$;
Create empty hash tables called $CompNbrs$ and $CompGraphs$;
Create a hash table called δ_{\max} and for each u in V set $\delta_{\max}(u)$ to zero;

```
1  $E_{\leq} \leftarrow$  array of edges of  $G$  that are sorted in order of increasing weight;  
2  $w_{old} \leftarrow$  weight of the lightest edge;  
3 for  $\{u, v\}$  in  $E_{\leq}$  do  
4    $w \leftarrow$  weight of  $\{u, v\}$ ;  
5   if  $w > w_{old}$  then  
6      $V_{flexible} \leftarrow Set()$ ;  
7     Add  $w$  to  $W$ ;  
8     for  $\{u, v\}$  in  $E_w$  do  
9       if  $Find(u) \neq Find(v)$  then  
10        Union( $u, v$ );  
11       else  
12        Add  $u$  to the set  $V_{flexible}$ ;  
13    $V_w \leftarrow$  vertices in  $E_w$ ;  
14   for  $u$  in  $V_w$  do  
15     Increase  $\delta_{\max}(u)$  by  $|CompNbrs(u)|$ ;  
16    $CompNbrs \leftarrow$  empty hash table;  
17   Add the set  $\{u, v, Comp(u, w), Comp(v, w)\}$  to the array  $CompGraphs(Find(u))$ , for each edge  
    $\{u, v\}$  in  $E_w$  such that  $Find(u)$  equals  $Find(x)$  for at least one vertex  $x$  in  $V_{flexible}$ ;  
18   Add to  $E_{fixed}$ , all the edges in  $E_w$  that are not in  $E_{flexible}$ ;  
19    $E_w \leftarrow$  empty array;  
20    $w_{old} \leftarrow w$ ;  
21   if  $Find(u) \neq Find(v)$  then  
22     Add  $\{u, v\}$  to  $E_w$ ;  
23     Set  $Comp(u, w)$  and  $Comp(v, w)$  to  $Find(u)$  and  $Find(v)$ , respectively;  
24     Add  $Comp(u, w)$  to the set  $CompNbr(v)$ ;  
25     Add  $Comp(v, w)$  to the set  $CompNbr(u)$ ;
```

26 If $|E_w|$ is greater than zero, then repeat lines 6 through 18;
27 Identify a ranking \mathcal{R} over V such that vertices with lower δ_{\max} are assigned higher ranks;
28 **for** $CompName$ in $CompGraphs.keys()$ **do**
29 $E_{comp} \leftarrow$ all edges $\{Comp(u, w), Comp(v, w)\}$ in $CompGraphs(CompName)$;
30 To each vertex $Comp(u, w)$ in E_{comp} , assign the rank $\mathcal{R}(u)$;
31 $E_{\leq comp} \leftarrow$ edges in E_{comp} sorted w.r.t. vertex rank (see Definition 2; each edge has weight w);
32 Add to $E_{selected}$, each edge in the graph that is constructed by applying Kruskal's algorithm to
 $E_{\leq comp}$;

Output: $M = (V, E_{fixed} \cup E_{selected})$

First we prove the correctness of Algo. 1, and subsequently, we derive its time complexity. Algo. 1 makes use of the disjoint-set data structure, which includes the operations: Make-Set, Find, and Union. The data structure is stored in memory in the form of a forest with self-loops and directed edges. Each directed edge from a vertex points to the parent of the vertex. A Make-Set operation creates a singleton vertex that points to itself. Each component in the forest has a single vertex that points to itself. This vertex is called the root. A Union operation takes as input, the roots of two components, and points one root to the

other. A Find operation takes as input a vertex, and returns the root of the component that contains the vertex. Specifically, we implemented balanced Union, and Find with path compression. For a more detailed description please read the survey by Galil and Italiano (1991).

Theorem 1. *Given as input a distance graph such that the distances are additive in some phylogenetic tree with strictly positive branch lengths, Algo. 1 constructs a vertex-ranked MST with the minimum number of leaves.*

Proof. Let $T = (V_T, E_T)$ be the phylogenetic tree that corresponds to the distance graph $G = (V, E)$. Let W be the set of weights of edges in E . Let \mathcal{F} be the laminar family over V , as defined in Lemma 3. Let g be the subgraph of G that contains the edges that are present in at least one MST of G . Let M be the output of Algo. 1.

Each edge in E_w is incident to vertices in different components. Since edges in E are visited in order of increasing weight, each edge in E_w is present in at least one MST of G .

Let c be the root of the component that is formed after Union operations are performed on each edge in E_w . Let E_c be the subset of E_w such that each edge in E_c is incident to vertices that are in component c after all Union operations on E_w have been performed. Let \mathcal{C} be the set of components such that each vertex in E_c is contained in a component in \mathcal{C} before any Union operations on E_w have been performed. Define the component graph $G_{\mathcal{C}}$ over \mathcal{C} to be the graph whose vertices are elements in \mathcal{C} , and whose edges are given by elements in E_c . It is clear that $G_{\mathcal{C}}$ is connected. We now consider the time point after all Union operations on E_w have been performed.

If $G_{\mathcal{C}}$ is a simple graph with no cycles, i.e., $|\mathcal{C}| = |E_c| - 1$, then each edge in E_c must be present in each MST of G . All edges in each simple, acyclic, component graph, are stored in E_{fixed} . If $G_{\mathcal{C}}$ is not simple, or if it contains cycles, then each edge in $G_{\mathcal{C}}$ is stored in $\text{CompGraphs}(c)$. Additionally each so-called component label $\{\text{Comp}(u, w), \text{Comp}(v, w)\}$ is also stored in $\text{CompGraphs}(c)$. For each vertex $u \in V$ the component label $\text{Comp}(u, w)$ is the root of the component that contains u before any union operations have been performed on edges in E_w . For each component c , the component graph $G_{\mathcal{C}}$ is induced by the component edges.

Let \mathcal{S} be the smallest sub-collection of the laminar family \mathcal{F} such that \mathcal{S} covers the neighbors of u but not u . Let F_w be the subgraph of G that is formed by removing from G all edges that are heavier than w . Let \mathcal{N}_w be the set of vertices in E_w that are adjacent to u . Let \mathcal{C}_w be the collection comprising the vertex set of each component of F_w that contains at least one vertex in \mathcal{N}_w . It is easy to see that $\mathcal{C}_w \subset \mathcal{S}$. It follows that $\mathcal{S} = \cup_{w \in W} \mathcal{C}_w$, where W is the set comprising the unique edge weights of G . Thus $\delta_{\max}(u) = |\mathcal{S}| = \sum_{w \in W} |\mathcal{C}_w|$. Thus the operations in line 15 correctly compute $\delta_{\max}(u)$.

At this time point all the edges of G have been visited. Subsequently, Algo. 1 selects a vertex ranking \mathcal{R} such that vertices with lower δ_{\max} are given higher ranks.

Let E_{flexible} be the set containing the edges $\{u, v\}$ that are stored in CompGraphs . Let Kruskal's algorithm be applied to the edges in $E_{\text{fixed}} \cup E_{\text{flexible}}$ that are sorted with respect to weight and \mathcal{R} , and let the resulting MST be the vertex-ranked MST $M_{\mathcal{R}} = (V_{\mathcal{R}}, E_{\mathcal{R}})$.

Let S be the set of all vertices in $\text{Comp}(u, w)$. From Lemma 2 (ii), we know that S induces a connected subgraph in each MST of G . This implies that, after all the edges that are no heavier than w have been visited by Algo. 1, the vertex set of the component that contains u is independent of the notion of the vertex rank that is used to sort the edges. Thus, instead of applying Kruskal's algorithm to each edge in $E_{\text{fixed}} \cup E_{\text{flexible}}$, we can avoid redundant computations by applying Kruskal's algorithm independently to each component graph. Consequently, $E_{\mathcal{R}} = E_{\text{fixed}} \cup E_{\text{selected}}$.

From Lemma 3 (ii), we know that, if there is a leaf l in $\mathbf{Sg}(h)$, such that $\{h, l\} \in E_T$, then among all vertices in $\mathbf{Sg}(h)$, $\delta_{\max}(l)$ is smallest. Consequently l has the highest rank in \mathcal{R} , when compared to other vertices in $\mathbf{Sg}(h)$. Since the surrogate vertex of h is the highest-ranked vertex in $\mathbf{Sg}(h)$, Algo. 1 implicitly selects l as the surrogate vertex of h . Since each leaf in T is adjacent to at most one hidden vertex, the vertex ranking that is selected by Algo. 1, maximizes the number of distinct leaves that are selected as surrogate vertices. Contracting the path in T between a hidden vertex and the corresponding surrogate vertex, increases the degree of the surrogate vertex. Thus, among all vertex-ranked MSTs, M has the minimum number of leaves. \square

7.4 Time complexity of Algorithm 1

We partition the operations of Algo. 1 into three parts. Part (i) sorts all the edges in E and performs Find and Union operations in order to select the edges in E_{fixed} and $CompGraphs$. Part (ii) computes δ_{\max} for each vertex in V , and part (iii) sorts, and applies Kruskal's algorithm to the edges in each component graph in $CompGraphs$.

In part (i) Algo. 1 iterates over the edges in G which are sorted w.r.t. edge weight. $G = (V, E)$ is a fully connected graph with n vertices and $n(n-1)/2$ edges. We used python's implementation of the Timsort algorithm (Peters, 2002) which sorts the edges in $O(n^2 \log n)$ time. Let m_f be the number of edges in E_{fixed} , and let m_c be the number of edges that are in a component graph. It is clear that $m_f + m_c \leq n(n-1)/2$. Algo. 1 iterates over each edge in G and performs $n(n-1)/2 + m_f + m_c$ Find operations, and $n-1$ Union operations. Since we implemented balanced Union, and Find with path compression, the time-complexity of these operations is $O((n(n-1)/2 + m_f + m_c)(\alpha((n(n-1)/2 + m_f + m_c, n))) = O(n^2(\alpha((n(n-1)/2 + m_f + m_c, n)))$, where $\alpha((n(n-1)/2 + m_f + m_c, n))$ is the inverse of Ackermann's function as defined in Tarjan (1975), and is less than 5 for all practical purposes. The total time complexity of part (i) is $O(n^2 \log n)$.

The operations in line 15 compute $\delta_{\max}(u)$ by counting the number of distinct components that cover the vertices $J_u \subset E_w$, such that each vertex $j \in J_u$ is adjacent to u . Assuming that the insertion and retrieval operations on hash tables, and insertion operations arrays have linear time-complexity, the total time complexity of part (ii) is $O(m_f + m_c)$.

Let the number of component graphs in $CompGraphs$ be k and let the number of edges and vertices in the i^{th} component graph be m_i and n_i , respectively. The time complexity of sorting, and applying Kruskal's algorithm to m_i edges, is $O(m_i \log m_i) + O(m_i \alpha(m_i, n_i)) = O(m_i \log m_i)$. The total time complexity of part (iii) is

$$\begin{aligned}
& \sum_{i=1}^k O(m_i \log m_i) \\
&= O\left(\sum_{i=1}^k m_i \log m_i\right) \\
&= O\left(\left(\sum_{i=1}^k m_i\right) \sum_{i=1}^k \frac{m_i}{\left(\sum_{i=1}^k m_i\right)} \log m_i\right) \\
&= O\left(m_c \sum_{i=1}^k \frac{m_i}{m_c} \log m_i\right) \\
&\leq O\left(m_c \log \sum_{i=1}^k \frac{m_i^2}{m_c}\right) \quad \text{from Jensen's inequality} \\
&\leq O\left(m_c \log \sum_{i=1}^k m_i^2\right) \\
&\leq O\left(m_c \log \left(\sum_{i=1}^k m_i\right)^2\right) \\
&= O(m_c \log m_c)
\end{aligned}$$

The total time complexity of Algo. 1 is $O(n^2 \log n) + O(m_f + m_c) + O(m_c \log m_c) = O(n^2 \log n)$.

8 Computational complexity of the MLVRMST construction problem

Let \mathcal{T} be the set of all phylogenetic trees. Let \mathcal{G} be the set of edge-weighted graphs, such that the edges of each graph in \mathcal{G} are weighted with distances that additive in some tree in \mathcal{T} . Algo. 1 constructs a MLVRMST of any graph in \mathcal{G} , in time $O(n^2 \log n)$. Thus, for graphs in \mathcal{G} , the decision version of the optimization problem MLVRMST is in the complexity class **P**. For graphs whose edges are not weighted with tree-additive distances, the MLVRMST problem may not be in **P**.

Consider the general optimization problem of constructing an MST with the minimum number of leaves (MLMST). Since the decision version of MLMST can be verified in polynomial time, MLMST is in **NP**. Additionally, it is easy to show that there is a polynomial time reduction from the Hamiltonian path problem to MLMST. Since the Hamiltonian path problem is in **NP-complete**, MLMST must be in **NP-hard** \cap **NP** = **NP-complete**.

9 Acknowledgements

We thank Erik Jan van Leeuwen and Davis Isaac for helpful discussions during the early stages of the work presented here.

10 Funding

PK's work has been funded in part by the German Center for Infection Research (DZIF, German Ministry of Education and Research Grants No. TTU 05.805, TTU 05.809).

11 Availability of code

A python implementation of Algo. 1 can be found at <http://resources.mpi-inf.mpg.de/departments/d3/publications/prabhavk/minLeavesVertexRankedMST>

References

- Buneman, P. 1971. The recovery of trees from measures of dissimilarity. In D. G. Kendall and P. Tautu, editors, *Mathematics in the Archaeological and Historical Sciences*, pages 387–395. Edinburgh University Press, Edinburgh, UK.
- Choi, M. J., Tan, V. Y. F., Anandkumar, A., and Willsky, A. S. 2011. Learning Latent Tree Graphical Models. *Journal of Machine Learning Research*, 12: 1771–1812.
- Chow, C. K. and Liu, C. N. 1968. Approximating discrete probability distributions with causal dependence trees. *IEEE Transactions on Information Theory*, IT-14(3): 462–467.
- dos Reis, M., Donoghue, P. C. J., and Yang, Z. 2016. Bayesian molecular clock dating of species divergences in the genomics era. *Nature Reviews Genetics*, 17(2): 71–80.
- Galil, Z. and Italiano, G. F. 1991. Data structures and algorithms for disjoint set union problems. *ACM Computing Surveys*, 23(3): 319–344.
- Gascuel, O. 1997. BIONJ: an improved version of the NJ algorithm based on a simple model of sequence data. *Molecular biology and evolution*, 14(7): 685–695.

- Huang, F., N., N. U., Perros, I., Chen, R., Sun, J., and Anandkumar, A. 2014. Scalable Latent Tree Model and its Application to Health Analytics. pages 1–19.
- Kalaghatgi, P., Pfeifer, N., and Lengauer, T. 2016. Family-joining: A fast distance-based method for constructing generally labeled trees. *Molecular Biology and Evolution*, 10(33): 2720–2734.
- Kruskal, J. B. 1956. On the Shortest Spanning Subtree of a Graph and the Traveling Salesman Problem. *Proceedings of the American Mathematical Society*, 7(1): 48–50.
- Peters, T. 2002. Timsort - Python. <https://svn.python.org/projects/python/trunk/Objects/listsort.txt>. See also <https://en.wikipedia.org/wiki/Timsort>.
- Ravi, R. and Singh, M. 2006. Delegate and conquer: An LP-based approximation algorithm for minimum degree MSTs. *Proceedings of the 33rd International Colloquium on Automata, Languages and Programming*, pages 169–180.
- Saitou, N. and Nei, M. 1987. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Molecular biology and evolution*, 4(4): 406–425.
- Tarjan, R. E. 1975. Efficiency of a Good But Not Linear Set Union Algorithm. *Journal of the ACM*, 22(2): 215–225.